



NLP RATIONALE

Adrián Salas Rodrigues



INTELLIGENT SYSTEMS

Problem to Solve	2
Experiment(s) done	6
Analysis of results	6

1. Problem to Solve

For this work I used a dataset of messages classified as spam or ham. What I did is to develop, using a BERT model and its features (text preprocessing and word-embedding), a spam classifier that is going to predict if a given message is spam or ham. The BERT model classifies the dataset by using pre-trained weights downloaded from the TensorFlow Hub repository.

	label	body
628	spam	New TEXTBUDDY Chat 2 horny guys in ur area 4 j...
1664	ham	Ãœ v ma fan...
1687	spam	Free Top ringtone -sub to weekly ringtone-get ...
1595	ham	Pls confirm the time to collect the cheque.
1263	ham	Ok. No wahala. Just remember that a friend in ...

The dataset can be found here:

<https://archive.ics.uci.edu/ml/datasets/SMS+Spam+Collection>

The GitHub code can be found here:

<https://github.com/AdrianSalas500/Spam-Classifier.git>

The steps for developing the spam classifier were:

1. Install dependencies

```
!pip install tensorflow
!pip install tensorflow-text
```

```
import tensorflow as tf
import tensorflow_hub as hub
import tensorflow_text as text
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
```

2. Import the dataset

```
df=pd.read_excel('SMSSpamCollection.xlsx')
```

3. Balance spam and ham classes and obtain a new balanced dataset

We have 747 spam emails and 4826 ham emails. The ham messages are significantly higher, implying that 15% are spam emails and 85% of ham emails, indicating an imbalance, so in order to balance the two classes, we reduce number of ham messages to 747.

```
df['label'].value_counts()
```

```
ham      4826
spam      747
Name: label, dtype: int64
```

3.1. Create two data frames (one for each class)

```
df_spam = df[df['label']=='spam']
df_ham = df[df['label']=='ham']
```

3.2. Balance the ham dataset and create a new balanced dataframe

```
df_ham_balanced = df_ham.sample(df_spam.shape[0])
```

```
df_balanced = pd.concat([df_ham_balanced, df_spam])
```

3.3. In this table you can see that now spam and ham have the same number of messages (747)

```
df_balanced['label'].value_counts()
```

```
spam      747
ham       747
Name: label, dtype: int64
```

4. Download the BERT models from the TensorFlow Hub repository (for preprocessing and encoding)

We download two BERT models, one to perform preprocessing and the other one for encoding.

```
bert_preprocess = hub.KerasLayer("https://tfhub.dev/tensorflow/bert_en_uncased_preprocess/3")
bert_encoder = hub.KerasLayer("https://tfhub.dev/tensorflow/bert_en_uncased_L-12_H-768_A-12/4")
```

5. Initialize the BERT and neural network layers

5.1. BERT layers

We use preprocess as the input for this layer. Then, the encoder is going to convert the preprocessed text in vectors (output of the layer).

```
text_input = tf.keras.layers.Input(shape=(), dtype=tf.string, name='text')
preprocessed_text = bert_preprocess(text_input)
outputs = bert_encoder(preprocessed_text)
```

5.2. Neural Network layers

The output is going to be fed in the neural network layers, that are two, the Dropout layer, and the Dense layer.

```
layer = tf.keras.layers.Dropout(0.1, name="dropout")(outputs['pooled_output'])
layer = tf.keras.layers.Dense(1, activation='sigmoid', name="output")(layer)
```

5.3. Final model

We add the input and output layers to construct the final model

```
model = tf.keras.Model(inputs=[text_input], outputs = [layer])
```

6. Split the dataset using the train_test_split function

Using sklearn, we split the dataset in two sets (training and testing sets)

```
X_train, X_test, y_train, y_test = train_test_split(df_balanced['body'], df_balanced['spam'], stratify=df_balanced['spam'])
model.fit(X_train, y_train, epochs=10)
```

7. Fit the model with 10 epochs

```
model.fit(X_train, y_train, epochs=10)
```

```

Epoch 1/10
35/35 [=====] - 187s 5s/step - loss: 0.5965 - accuracy: 0.7214
Epoch 2/10
35/35 [=====] - 218s 6s/step - loss: 0.4863 - accuracy: 0.8170
Epoch 3/10
35/35 [=====] - 244s 7s/step - loss: 0.4209 - accuracy: 0.8589
Epoch 4/10
35/35 [=====] - 228s 7s/step - loss: 0.3780 - accuracy: 0.8687
Epoch 5/10
35/35 [=====] - 233s 7s/step - loss: 0.3421 - accuracy: 0.8902
Epoch 6/10
35/35 [=====] - 228s 7s/step - loss: 0.3192 - accuracy: 0.8938
Epoch 7/10
35/35 [=====] - 231s 7s/step - loss: 0.3089 - accuracy: 0.9027
Epoch 8/10
35/35 [=====] - 242s 7s/step - loss: 0.2848 - accuracy: 0.9089
Epoch 9/10
35/35 [=====] - 230s 7s/step - loss: 0.2736 - accuracy: 0.9098
Epoch 10/10
35/35 [=====] - 225s 6s/step - loss: 0.2638 - accuracy: 0.9134

```

The model obtains a **91,34 %** of accuracy.

8. Evaluate the model in the testing dataset to obtain an array of 0's and 1's predicting if a message is spam or ham

```

y_pred = model.predict(X_test)
y_pred = y_pred.flatten()

```

```

y_pred = np.where(y_pred > 0.5, 1, 0)
y_pred

```

```

array([0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0,
       1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1,
       0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1,
       1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 1,
       1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0,
       0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0,
       1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1,
       1, 1, 0, 0, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0,
       1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1,
       0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 1,
       0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0,
       0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0,
       1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1,
       0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1,
       0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0,
       1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0])

```

2. Experiment(s) done

To assess that the model classified correctly, test messages were used.

```
sample_dataset = [  
    'You can win a lot of money, register in the link below',  
    'You have an iPhone 10, spin the image below to claim your prize and it will be delivered in your door step',  
    'You have an offer, the company will give you 50% off on every item purchased.',  
    'Hey Bravin, do not be late for the meeting tomorrow will start lot exactly 10:30 am',  
    "See you monday, we have alot to talk about the future of this company ."  
]
```

This array of samples is composed by three spam messages and two ham ones.

3. Analysis of results

After defining the samples, I used the model to predict the results. These results were an array of numbers, in which a number above 0.5 indicated that the message is considered spam, and a number below 0.5 that is ham.

```
array([[0.69952023],  
       [0.9169347 ],  
       [0.5298996 ],  
       [0.71341854],  
       [0.15887296]], dtype=float32)
```

As can be seen in the image, 4 out of 5 samples were correctly predicted, giving an accuracy rate of 80 %.